

Notes on Matrix Derivative and Trace

Wu Xiaojian

July 27, 2019

1 Notation

Let

$$f : \mathbb{R}^{m \times n} \mapsto \mathbb{R} \quad (1)$$

be a scalar function, taking matrix of $m \times n$ shape as input, and throwing out a real number as output. By convention, we use bold capitalized letters such as \mathbf{X}, \mathbf{Y} to denote matrix, as follows

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,n} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \cdots & x_{m,n} \end{bmatrix}$$

We define the derivative of f with respect to \mathbf{X} as

$$\frac{\partial f}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial f}{\partial x_{1,1}} & \frac{\partial f}{\partial x_{1,2}} & \cdots & \frac{\partial f}{\partial x_{1,n}} \\ \frac{\partial f}{\partial x_{2,1}} & \frac{\partial f}{\partial x_{2,2}} & \cdots & \frac{\partial f}{\partial x_{2,n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{m,1}} & \frac{\partial f}{\partial x_{m,2}} & \cdots & \frac{\partial f}{\partial x_{m,n}} \end{bmatrix} \quad (2)$$

that is, the derivative applies to each elements of \mathbf{X} and is still a matrix with the same shape as \mathbf{X} , much like to the broadcast operation upon a matrix which applies the specified operation to all the elements and remains the shape. Similarly, we define the differential of \mathbf{X} as

$$d\mathbf{X} = \begin{bmatrix} dx_{1,1} & dx_{1,2} & \cdots & dx_{1,n} \\ dx_{2,1} & dx_{2,2} & \cdots & dx_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ dx_{m,1} & dx_{m,2} & \cdots & dx_{m,n} \end{bmatrix} \quad (3)$$

Finally, let's introduce *trace*. For a square $\mathbf{X}^{n \times n}$, we use $\text{tr}(\mathbf{X})$ to represent the trace of \mathbf{X} , the sum of those among in the main diagonal of \mathbf{X}

$$\text{tr}(\mathbf{X}) = \sum_{i=1}^n x_{i,i} \quad (4)$$

which shows¹ $\text{tr}(\mathbf{X}) \in \mathbb{R}$.

2 Inner Product of Matrix

For any $\mathbf{v}, \mathbf{k} \in \mathbb{R}^n$, their inner product can be defined as

$$\langle \mathbf{v}, \mathbf{k} \rangle = \mathbf{v}^T \mathbf{k} = \sum_{i=1}^n v_i k_i$$

It's very natural to borrow this concept and define the *inner product of matrix*. Suppose \mathbf{A}, \mathbf{B} are matrix of the same shape $m \times n$, their inner product is

$$\langle \mathbf{A}, \mathbf{B} \rangle_{m \times n} = \sum_{\substack{1 \leq j \leq m \\ 1 \leq i \leq n}} a_{i,j} b_{i,j}$$

that is, the sum of product entrywise of two matrix. Now that we have defined trace, we can use trace to define inner product of matrix, in a elegant way.

Theorem 2.1 Suppose \mathbf{A}, \mathbf{B} are matrix of the same shape $m \times n$, then

$$\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^T \mathbf{B}) \quad (5) \quad \square$$

Proof Suppose $\mathbf{A} = [\alpha_1, \alpha_2, \dots, \alpha_n]$, $\mathbf{B} = [\beta_1, \beta_2, \dots, \beta_n]$ where α_i, β_i are column vectors, then $\mathbf{A}^T = [\alpha_1^T, \alpha_2^T, \dots, \alpha_n^T]^T$, and

$$\mathbf{A}^T \mathbf{B} = \begin{bmatrix} \alpha_1^T \\ \alpha_2^T \\ \vdots \\ \alpha_n^T \end{bmatrix} [\beta_1, \beta_2, \dots, \beta_n] = \begin{bmatrix} \alpha_1^T \beta_1 & & & \\ & \alpha_2^T \beta_2 & & \\ & & \ddots & \\ & & & \alpha_n^T \beta_n \end{bmatrix}$$

the other entries of which are not important and can be omitted. Then

$$\text{tr}(\mathbf{A}^T \mathbf{B}) = \sum_{j=1}^n \alpha_j^T \beta_j = \sum_{j=1}^n \sum_{i=1}^m a_{ij} b_{ij} = \langle \mathbf{A}, \mathbf{B} \rangle$$

3 Derivative of Multi-variable Function

Recall that in multi-variable calculus, the *complete differential* of f with respect to variables x_1, x_2, \dots, x_n is

$$df = \sum_{i=1}^n \frac{\partial f}{\partial x_i} dx_i$$

¹We assume all matrix we discussing here are real matrix.

Using nabla operator $\nabla = [\frac{\partial}{\partial x_1}, \frac{\partial}{\partial x_2}, \dots, \frac{\partial}{\partial x_n}]^T$, we rewrite the complete differential above as more gracefully as

$$df = \nabla f \cdot d\mathbf{x} = (\nabla f)^T d\mathbf{x} \quad (6)$$

where

$$\nabla f = [\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n}]^T \quad \text{and} \quad d\mathbf{x} = [dx_1, dx_2, \dots, dx_n]^T$$

Now we can generalize this concept to matrix, say, the complete differential of f with respect to $\mathbf{X}_{m \times n}$, and so we have

$$df|_{\mathbf{X}} = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial x_{i,j}} dx_{i,j} \quad (7)$$

in most cases, if context agrees, we can denote it shortly as df .

Example 3.1 Suppose

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \end{bmatrix}$$

we have

$$\frac{\partial f}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \frac{\partial f}{\partial x_{12}} & \frac{\partial f}{\partial x_{13}} \\ \frac{\partial f}{\partial x_{21}} & \frac{\partial f}{\partial x_{22}} & \frac{\partial f}{\partial x_{23}} \end{bmatrix} \quad \text{and} \quad d\mathbf{X} = \begin{bmatrix} dx_{11} & dx_{12} & dx_{13} \\ dx_{21} & dx_{22} & dx_{23} \end{bmatrix}$$

Let's transpose the derivative and multiply them, and we can get

$$\left(\frac{\partial f}{\partial \mathbf{X}} \right)^T d\mathbf{X} = \begin{bmatrix} \frac{\partial f}{\partial x_{11}} dx_{11} + \frac{\partial f}{\partial x_{21}} dx_{21} & \frac{\partial f}{\partial x_{12}} dx_{12} + \frac{\partial f}{\partial x_{22}} dx_{22} & \frac{\partial f}{\partial x_{13}} dx_{13} + \frac{\partial f}{\partial x_{23}} dx_{23} \end{bmatrix}$$

It's very obvious that

$$\text{tr} \left(\left(\frac{\partial f}{\partial \mathbf{X}} \right)^T d\mathbf{X} \right) = \sum_{i=1}^2 \sum_{j=1}^3 \frac{\partial f}{\partial x_{ij}} dx_{ij} = df \quad (8)$$

In EXAMPLE(3.1) we find it seems that the trace of the multiplication between the transpose of derivative of f with respect to \mathbf{X} and the differential of \mathbf{X} equals to the differential of f . In fact, we have the theorem, as follows

Theorem 3.1 Suppose $\mathbf{X} \in \mathbb{R}^{m \times n}$ and $f : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$, then

$$df = \langle \frac{\partial f}{\partial \mathbf{X}}, d\mathbf{X} \rangle = \text{tr} \left(\left(\frac{\partial f}{\partial \mathbf{X}} \right)^T d\mathbf{X} \right) \quad (9) \quad \square$$

Proof By THEOREM(2.1), we have

$$\langle \frac{\partial f}{\partial \mathbf{X}}, d\mathbf{X} \rangle = \text{tr} \left(\left(\frac{\partial f}{\partial \mathbf{X}} \right)^T d\mathbf{X} \right)$$

On the other hand, since the inner product of two matrix is the sum of each product entrywise of the two, that is,

$$\langle \frac{\partial f}{\partial \mathbf{X}}, d\mathbf{X} \rangle = \sum_{i,j} \frac{\partial f}{\partial x_{i,j}} dx_{i,j}$$

we have $df = \langle \frac{\partial f}{\partial \mathbf{X}}, d\mathbf{X} \rangle$, and thus the theorem holds. ■

4 Matrix Differential

Analogous to laws of derivative on real-valued function, we can define the derivative of matrix and derive some operation laws from it.

Theorem 4.1 (Basic Laws of Arithmetic Operations)

$$d(\mathbf{X} \pm \mathbf{Y}) = d\mathbf{X} \pm d\mathbf{Y} \quad (10)$$

$$d(\mathbf{XY}) = (d\mathbf{X})\mathbf{Y} + \mathbf{X}(d\mathbf{Y}) \quad (11)$$

$$d(\mathbf{X} \odot \mathbf{Y}) = d(\mathbf{X}) \odot \mathbf{Y} + \mathbf{X} \odot d(\mathbf{Y}) \quad (12)$$

where \odot is **Hadamard Product**(also known as the **Schur product** or the **Entrywise Product**) is a binary operation that takes two matrices of the same dimensions, and produces another matrix where each element ij is the product of elements ij of the original two matrices. It should not be confused with the more common matrix product. □

Proof

$$\begin{aligned} d(\mathbf{X} \pm \mathbf{Y}) &= \begin{bmatrix} d(x_{11} \pm y_{11}) & \cdots & d(x_{1n} \pm y_{1n}) \\ \vdots & \ddots & \vdots \\ d(x_{m1} \pm y_{m1}) & \cdots & d(x_{mn} \pm y_{mn}) \end{bmatrix} \\ &= \begin{bmatrix} dx_{11} \pm dy_{11} & \cdots & dx_{1n} \pm dy_{1n} \\ \vdots & \ddots & \vdots \\ dx_{m1} \pm dy_{m1} & \cdots & dx_{mn} \pm dy_{mn} \end{bmatrix} \\ &= \begin{bmatrix} dx_{11} & \cdots & dx_{1n} \\ \vdots & \ddots & \vdots \\ dx_{m1} & \cdots & dx_{mn} \end{bmatrix} \pm \begin{bmatrix} dy_{11} & \cdots & dy_{1n} \\ \vdots & \ddots & \vdots \\ dy_{m1} & \cdots & dy_{mn} \end{bmatrix} \\ &= d\mathbf{X} \pm d\mathbf{Y} \end{aligned}$$

$$d(\mathbf{XY}) = d \left(\begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ \vdots & \vdots & \cdots & \vdots \\ x_{m1} & x_{n2} & \cdots & x_{mn} \end{bmatrix} \begin{bmatrix} y_{11} & \cdots & y_{1t} \\ y_{21} & \cdots & y_{2t} \\ \vdots & \cdots & \vdots \\ y_{n1} & \cdots & y_{nt} \end{bmatrix} \right) = d \begin{bmatrix} z_{11} & \cdots & z_{1t} \\ \vdots & \ddots & \vdots \\ z_{n1} & \cdots & z_{nt} \end{bmatrix}$$

where $z_{ij} = \sum_{k=1}^n x_{ik}y_{kj}$, so

$$dz_{ij} = d \sum_{k=1}^n x_{ik}y_{kj} = \sum_{k=1}^n d(x_{ik})y_{kj} + \sum_{k=1}^n x_{ik}d(y_{kj})$$

For the last two items, we can convert each of them back to the form of matrix multiplication and so we will get

$$d\mathbf{Z} = d(\mathbf{X})\mathbf{Y} + \mathbf{X}d(\mathbf{Y})$$

The proof of the hadamard product is similar. ■

Theorem 4.2

$$d(\mathbf{X}^T) = (d\mathbf{X})^T \quad (13)$$

$$d(\text{tr}(\mathbf{X})) = \text{tr}(d\mathbf{X}) \quad (14)$$

Theorem 4.3 (Differential of Constant Matrix) Assume \mathbf{C} is a constant matrix, then

$$d(\mathbf{C}) = \mathbf{O} \quad (15) \quad \square$$

Theorem 4.4 (Differential of Inverse of Matrix) Suppose \mathbf{X} is invertible, and its inverse is \mathbf{X}^{-1} , then

$$d(\mathbf{X}^{-1}) = -\mathbf{X}^{-1}d(\mathbf{X})\mathbf{X}^{-1} \quad (16) \quad \square$$

Proof Let \mathbf{I} be identity. Since \mathbf{X} is invertible, therefore $d\mathbf{I} = d(\mathbf{X}\mathbf{X}^{-1})$. Notice \mathbf{I} is a constant matrix, thus $d\mathbf{I} = \mathbf{O}$ and

$$d(\mathbf{X}\mathbf{X}^{-1}) = d(\mathbf{X})\mathbf{X}^{-1} + \mathbf{X}d(\mathbf{X}^{-1}) = \mathbf{O}$$

This concludes that $d(\mathbf{X}^{-1}) = -\mathbf{X}^{-1}d(\mathbf{X})\mathbf{X}^{-1}$. ■

Theorem 4.5 (Chain Rule) Let σ is a scalar function, and we define the broadcast operation as

$$\sigma(\mathbf{X}) = \begin{bmatrix} \sigma(x_{11}) & \cdots & \sigma(x_{1n}) \\ \vdots & \ddots & \vdots \\ \sigma(x_{m1}) & \cdots & \sigma(x_{mn}) \end{bmatrix}$$

then

$$d(\sigma(\mathbf{X})) = \sigma'(\mathbf{X}) \odot d\mathbf{X} \quad (17) \quad \square$$

Proof

$$\begin{aligned}
d(\sigma(\mathbf{X})) &= \begin{bmatrix} d\sigma(x_{11}) & \cdots & d\sigma(x_{1n}) \\ \vdots & \ddots & \vdots \\ d\sigma(x_{m1}) & \cdots & d\sigma(x_{mn}) \end{bmatrix} \\
&= \begin{bmatrix} \sigma'(x_{11})dx_{11} & \cdots & \sigma'(x_{1n})dx_{1n} \\ \vdots & \ddots & \vdots \\ \sigma'(x_{m1})dx_{m1} & \cdots & \sigma'(x_{mn})dx_{mn} \end{bmatrix} \\
&= \begin{bmatrix} \sigma'(x_{11}) & \cdots & \sigma'(x_{1n}) \\ \vdots & \ddots & \vdots \\ \sigma'(x_{m1}) & \cdots & \sigma'(x_{mn}) \end{bmatrix} \odot \begin{bmatrix} dx_{11} & \cdots & dx_{1n} \\ \vdots & \ddots & \vdots \\ dx_{m1} & \cdots & dx_{mn} \end{bmatrix} \\
&= \sigma'(\mathbf{X}) \odot d\mathbf{X}
\end{aligned}$$

5 Trace

Theorem 5.1 Suppose a, b, c is scalar value, $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{m \times n}$, $\mathbf{Z} \in \mathbb{R}^{n \times m}$. Here are some properties about trace

$$\text{tr}(c) = c \quad (18)$$

$$\text{tr}(\mathbf{X}^T) = \text{tr}(\mathbf{X}) \quad (19)$$

$$\text{tr}(a\mathbf{X} \pm b\mathbf{Y}) = a\text{tr}(\mathbf{X}) \pm b\text{tr}(\mathbf{Y}) \quad (20)$$

$$\text{tr}(\mathbf{XZ}) = \text{tr}(\mathbf{ZX}) \quad (21)$$

Proof we just proof the last one. Since the (i, j) -entry of \mathbf{XZ} is $\sum_{k=1}^n x_{ik}z_{kj}$, and that of \mathbf{ZX} is $\sum_{p=1}^m z_{ip}x_{pj}$, thus

$$\text{tr}(\mathbf{XZ}) = \sum_{p=1}^m \sum_{k=1}^n x_{pk}z_{kp} = \sum_{k=1}^n \sum_{p=1}^m z_{kp}x_{pk} = \text{tr}(\mathbf{ZX})$$

We can also infer from this property that if $\mathbf{X}_1\mathbf{X}_2 \cdots \mathbf{X}_n$ and $\mathbf{X}_n\mathbf{X}_1 \cdots \mathbf{X}_{n-1}$ are both defined, then

$$\text{tr}(\mathbf{X}_1\mathbf{X}_2 \cdots \mathbf{X}_n) = \text{tr}(\mathbf{X}_n\mathbf{X}_1 \cdots \mathbf{X}_{n-1}) \quad (22)$$

The proof is very easy that using $\text{tr}(\mathbf{XZ}) = \text{tr}(\mathbf{ZX})$, let $\mathbf{X} = \mathbf{X}_1\mathbf{X}_2 \cdots \mathbf{X}_{n-1}$ and $\mathbf{Z} = \mathbf{X}_n$ and so we make it. \blacksquare